# Verification of land-atmosphere coupling in forecast models, reanalyses and land surface models using flux site observations

Paul A. Dirmeyer[1]*, Liang Chen[1], Jiexia Wu[1], Chul-Su Shin[1], Bohua Huang[1], Benjamin A. Cash[1], Michael G. Bosilovich[2], Sarith Mahanama[2], Randal D. Koster[2], Joseph A. Santanello[2], Michael B. Ek[3], Gianpaolo Balsamo[4], Emanuel Dutra[5], and D. M. Lawrence[6]

[1]Center for Ocean-Land-Atmosphere Studies, George Mason University

[2]NASA / Goddard Space Flight Center

[3]NOAA / National Centers for Environmental Prediction / Environmental Modeling Center

[4]European Centre for Medium-range Weather Forecasts

[5]Instituto Dom Luiz, Faculdade de Ciências, Universidade de Lisboa[5]National Center for Atmospheric Research

*Corresponding Author:

Paul A. Dirmeyer
Center for Ocean-Land-Atmosphere Studies
George Mason University
4400 University Drive, Mail Stop: 6C5
Fairfax, Virginia 22030  USA
*pdirmeye@gmu.edu*

Submitted to:  *Journal of Hydrometeorology*

1  **Abstract:**

2     We confront four model systems in three configurations (LSM, LSM+GCM, and

3  reanalysis) with global flux tower observations to validate states, surface fluxes, and

4  coupling indices between land and atmosphere. Models clearly under-represent the

5  feedback of surface fluxes on boundary layer properties (the atmospheric leg of land-

6  atmosphere coupling), and may over-represent the connection between soil moisture and

7  surface fluxes (the terrestrial leg). Models generally under-represent spatial and temporal

8  variability relative to observations, which is at least partially an artifact of the differences

9  in spatial scale between model grid boxes and flux tower footprints. All models bias high in

10  near-surface humidity and downward shortwave radiation, struggle to represent

11  precipitation accurately, and show serious problems in reproducing surface albedos. These

12  errors create challenges for models to partition surface energy properly and errors are

13  traceable through the surface energy and water cycles. The spatial distribution of the

14  amplitude and phase of annual cycles (first harmonic) are generally well reproduced, but

15  the biases in means tend to reflect in these amplitudes. Interannual variability is also a

16  challenge for models to reproduce. Our analysis illuminates targets for coupled land-

17  atmosphere model development, as well as the value of long-term globally-distributed

18  observational monitoring.

19

## 1. Introduction

Many LSMs were developed and pressed into service during the 1980s and 1990s to provide lower boundary conditions for the atmospheric GCMs used in climate and weather simulation and prediction (Santanello et al. 2017). This occurred at a time when observations of key land surface variables, and the coupled processes that link the water and energy cycles between the land and atmosphere, were extremely limited. As a result, performance of coupled LSM-GCM systems has been sub-optimal (Dirmeyer et al. 2017).

The necessary observational data sets for validation are only recently becoming available; datasets that combine co-located measurements of land surface states, surface fluxes, near-surface meteorology, and properties of the atmospheric column. Early field campaigns (e.g., Sellers et al. 1992, 1995; Famiglietti et al. 1999; Jackson and Hsu 2001; Andreae 2002) provided observations that helped advance theory and model parameterization development, but their short periods of operation meant collected data provided limited sampling of the phase-space of land-atmosphere interactions, rarely quantifying interannual variability. In the mid-1990s, networks of observing stations began to be established and maintained, providing long-term data sets. A growing number of soil moisture monitoring networks have been established. Their data have been collated, homogenized and standardized by two separate efforts (Dorigo et al. 2011, 2013, 2017; Quiring et al. 2016). Those data sets were used by Dirmeyer et al. (2016) in a first-of-its-kind multi-model multi-configuration assessment of soil moisture simulation fidelity.

Simultaneously, efforts began in the ecological community to collect surface flux data over a variety of biomes (FLUXNET; Baldocchi et al 2001). Over time, in consultation with interested scientific communities, FLUXNET expanded their instrumentation suite to measure soil moisture, ground heat flux, and four-component radiation, allowing detailed

closure of the surface energy balance. Rigid standards for data formatting and dissemination within and across regional networks was lacking, so a global standardized and quality-controlled subset of data from many FLUXNET sites was produced ("La Thuile FLUXNET dataset", cf. http://www.fluxdata.org) covering multiple links in the coupled land-atmosphere process chain (Santanello et al. 2011). The La Thuile data set enabled a greater degree of model validation (e.g., Williams et al. 2009; Bonan et al. 2012; Boussetta et al. 2013; Melaas et al. 2013; Balzarolo et al. 2014; Purdy et al. 2016).

In this study, we employ the updated FLUXNET2015 synthesis data set, (Pastorello et al. 2017) expanding the multi-model multi-configuration study of soil moisture simulations in Dirmeyer et al. (2016) to a global assessment of surface energy and water balance simulations, and basic metrics of land-atmosphere coupling. Section 2 describes the observational data and models examined. The next three sections present validations of model annual means, annual cycles, and coupling metrics. We then discuss some of the pathological model behaviors that emerge from the analysis and present conclusions. Throughout the paper we present synthesis figures. Detailed scatter plots showing results across all FLUXNET2015 sites for each model are consigned to the Supplement.

## 2. Data and Models

The range of dates of data varies considerably among model simulations, and also between individual observational sites. We analyze spatial variability and compare only climatologies (annual means or mean annual cycles) in order to minimize the effect of such asynchronicities, and present a quantification of interannual variability. It is not the intent of this study to validate model simulations of specific events, but rather their overall coupled land-atmosphere behavior. Note also that many coupling metrics, including those

68 used here, can be calculated for LSMs from a combination of forcing and model output,

69 even though the LSMs are not coupled to GCMs.

70 *2.1 Observed data*

71 In situ measurements of near surface meteorological variables, surface fluxes and soil

72 moisture used for model validation come from the November 2016 version of the

73 FLUXNET2015 station data set. Daily, monthly and yearly data have been used; processing

74 of the meteorological, radiation, heat flux and surface hydrologic data including gap-filling

75 are described by Reichstein et al. (2005) and Vuichard and Papale (2015). Only the Tier 1

76 (open access) data are used in this study (see Table S1 for a complete list of sites) – Figure

77 1 shows the spatial distribution of sites and some of the key characteristics regarding data

78 availability. 166 sites provide 1242 site-years of data, but coverage is concentrated in the

79 mid-latitudes and particular underrepresentation in the tropics.

80 The variables processed for this analysis include surface pressure, near surface air

81 temperature and vapor pressure deficit, precipitation, four-component and net radiation,

82 surface sensible and latent heat fluxes (gap-filled following the method of Reichstein et al.

83 2005 and energy balance closure-corrected) and soil water content measured at the first

84 (shallowest) sensor. There is no consolidated information on the depth of the shallowest

85 sensor across all sites, but typically it is at 5cm or 10cm below the surface. Vapor pressure

86 deficit is converted to specific humidity using the Clausius-Clapeyron relationship. We have

87 used the provided FLUXNET2015 data at the corresponding time intervals for each

88 calculation: yearly data for annual means, monthly data for annual cycles, and daily data for

89 calculating coupling indices.

90 In addition, we examine a number of gridded global precipitation products for

91 comparison to FLUXNET2015 sites. These are listed in Table S2.

92    *2.2 Model systems*

93    Four global modeling systems are evaluated; two from operational forecast centers and

94    two that are primarily used for research. The operational systems are from the U.S.

95    National Oceanic and Atmospheric Administration (NOAA) National Centers for

96    Environmental Prediction (NCEP) and the European Centre for Medium-range Weather

97    Forecasts (ECMWF). The research systems are from the U.S. National Aeronautics and

98    Space Administration (NASA) Global Modeling and Assimilation Office (GMAO) and the U.S.

99    National Center for Atmospheric Research (NCAR).

100   Table 1 summarizes the model components and configurations. Generally, each

101   modeling system is interrogated in three different configurations: 1) LSM only (offline),

102   driven by gridded observationally-based meteorological analyses including downward

103   radiation; 2) LSM coupled to GCM in a *free-running* mode where the coupled system

104   evolves unconstrained after initialization; 3) Reanalysis, where the coupled LSM and GCM

105   are constrained by data assimilation at diurnal or sub-diurnal increments to represent the

106   actual historical evolution of state variables. The NCAR model system does not have an

107   associated reanalysis, so to keep the four-by-three matrix filled, two different reanalyses

108   from GMAO are included. Note that when the coordinates for a FLUXNET2015 site lie

109   within a model's ocean grid cell, it is excluded from comparisons for that model. Thus, the

110   number of stations compared vary from model to model depending on resolution and the

111   land-sea mask.

112   2.2.1 NCEP

113   Data for the offline configuration comes from an author-produced simulation using Noah

114   LSM version 2.7.1 (Ek et al., 2003, Mitchell, 2005) driven by 3-hourly gridded

115   meteorological data from the Terrestrial Hydrology Research Group at Princeton

4

116    University (Sheffield et al., 2006). The free-running coupled land-atmosphere simulation

117    consists of a subset of 48 years from a 420 year long current climate simulation of CFSv2

118    initialized in 1980 (Shukla et al. 2017). The coupled simulation is unique among the model

119    systems in that it also includes a coupled ocean component. However, this should have very

120    little effect on the local coupled land-atmosphere behavior of the model. Years 2101-2148

121    of the simulation are used, but the calendar dates have no real meaning in a fully coupled

122    climate model so far from the initial state, wherein attributes such as atmospheric

123    composition, solar intensity, orbital parameters, etc., are held constant at late 20th century

124    values. The latest NCEP reanalysis is also examined (CFSR; Saha et al. 2010), which

125    combines a global land data assimilation system derived from the NASA Land Information

126    System (LIS; Peters-Lidard et al., 2007), driven by a blended global precipitation analysis

127    (Xie and Arkin 1997; Xie et al. 2007), used to update the coupled analysis cycle once per

128    day over the period 1979-2009.

129    2.2.2 GMAO

130    Two reanalyses are included for GMAO; version 1 and version 2 of the Modern-Era

131    Retrospective Analysis for Research and Applications (MERRA; Rienecker et al. 2011,

132    Reichle et al. 2017a). MERRA data cover the period 1980-2015. MERRA-2 is the current

133    state-of-the-art reanalysis covering 1980-2015 (Molod et al. 2015, Gelaro et al. 2017), and

134    is the source of most of the meteorological forcing data for the offline simulation of the

135    Catchment LSM v25 C05 (GMAO 2015a,b). As part of the MERRA-2 reanalysis, the GCM-

136    generated precipitation is corrected with observations-based precipitation before it

137    reaches the land surface (Reichle et al. 2017b); the reanalysis meteorological fields thus

138    feel the observed precipitation rates indirectly through the surface fluxes. Additionally, a

139    global 36-year offline Catchment simulation on the MERRA grid and a 16-year coupled

140 GEOS5-Catchment simulation at half-degree resolution with prescribed observed SSTs

141 were generated for this comparison.

142 2.2.3 NCAR

143 There is no operational reanalysis produced with the NCAR Community Earth System

144 Model (CESM). However, CESM is widely used for research in the academic community,

145 and we have generated offline and coupled simulations for this comparison. The offline

146 simulation uses version 4.5 of the Community Land Model (CLM; Lawrence et al. 2011)

147 driven with forcing spanning 1991-2010 from version 4 of the blended and gap-filled

148 CRUNCEP (Viovy 2013) 0.5° data set (available at:

149 https://www.earthsystemgrid.org/dataset/ucar.cgd.ccsm4.CRUNCEP.v4.html) aggregated

150 to the nominal 1° GCM resolution. A simulation with CLM4.5 coupled to CAM4 in CESM1.2.2

151 has been produced spanning 1991-2014 with specified climatological SSTs.

152 2.2.4 ECMWF

153 The offline simulation from ECMWF is with Cycle 43R1 of the Hydrology Tiled ECMWF

154 Scheme of Surface Exchanges over Land (HTESSEL) run at ~16km resolution based on a

155 cubic octahedral global grid (TCo639) for the period 1979-2015. This offline simulation

156 follows ERA-Interim/land configurations closely (see Balsamo et al. 2015), forced by ERA-

157 Interim meteorology and fluxes with an altitude correction applied to temperature,

158 humidity and surface pressure. This offline simulation is used to initialized the land state of

159 the operational ECMWF hindcasts. The coupled simulation comes from the Athena Project

160 (Kinter et al. 2013) for 1961-2007 where an older version of HTESSEL is coupled to IFS

161 Cycle 32R3 at a similarly high native horizontal resolution and specified observed SSTs, but

162 the data has been post-processed to a 1.125° uniform grid. ERA-Interim (Dee et al. 2011),

163    spanning 1979-2015, provides the reanalysis configuration of data for the comparison,

164    which used TESSEL prior to hydrology upgrades.

165

166    **3. Annual Means**

167    The comparison of models to FLUXNET2015 observations of annual means amounts to

168    an assessment of model ability to reproduce global spatial patterns (within the limitations

169    of the uneven distribution of station locations) of the variables' time averages. For the

170    offline LSM simulations, meteorological forcing data are specified from gridded data sets,

171    so their correlation to FLUXNET2015 observations is not a pure reflection of model

172    performance as the forcing data constrain LSM behavior. Similarly, for the reanalysis

173    products, performance reflects a combination of model characteristics, data assimilation

174    techniques and the distribution and quality of the data assimilated. Assimilation of

175    observational data constrains the coupled land-atmosphere model behavior to some

176    degree. While the free-running model simulations provide an unabridged assessment of

177    model performance, results from the other modes of simulation are nevertheless

178    enlightening.

179    As an indicator of observational uncertainty and the impact of comparing model grid

180    box values to field sites, we first note how a number of gridded observational precipitation

181    products and the reanalyses validate against precipitation measurements at FLUXNET2015

182    locations. Figure 2 shows mean (dots) and span (whiskers) of annual precipitation totals,

183    where the abscissa always corresponds to measurements from the FLUXNET2015 sites.

184    For most sites, the observational products (top two rows of Fig. 2) cover the entire time

185    span of FLUXNET2015 observations (see Table S2 for details). All reanalyses (bottom row

186    of Fig. 2) except CFSR span the FLUXNET2015 period. Several statistics of spatial

7

187     agreement are shown: Pearson's product moment correlation coefficient ($r_p$), Spearman's

188     rank correlation coefficient ($r_s$), root mean square error (RMSE), slope of the best-fit linear

189     regression of Y on X (Slope) and the fraction of total stations (labeled "Span Diag" in Fig. 2)

190     where the span of the individual annual totals from the gridded products (vertical

191     whiskers) overlap the span from FLUXNET2015 sites (horizontal whiskers). The last

192     statistic tests the possibility that the FLUXNET2015 observations and gridded estimates do

193     not come from distinct populations, i.e. their ranges overlap.

194     Estimates from gridded observational data sets, which range in spatial resolution from

195     0.25° (MSWEP, TRMM) to 2.5° (GPCP), provide a plausible upper bound to the accuracy we

196     could expect from gridded Earth system models. For the 166 (or fewer) FLUXNET2015

197     sites compared, which admittedly represent a rather uneven sampling of global terrestrial

198     precipitation, three observational products score at the top: MSWEP, CPC-Uni and U.Del.

199     Each has a Pearson's correlation of nearly 0.8, a rank correlation between 0.8-0.9, and the

200     highest number of stations whose ranges span the diagonal X=Y line. The lower limit for

201     RMSE across these sites is about 240mm. Note that all gridded products underestimate the

202     slope, indicating the inability of large area averages to resolve local variations in average

203     precipitation.

204     MERRA-2 performs on par with the best gridded observed products, namely because it

205     reports a bias corrected precipitation that is used as part of the assimilation process

206     instead of model-generated precipitation as an input to the LSM (Reichle and Liu 2014).

207     Thus, it is effectively another gridded observational data set for precipitation. Figure S1

208     compares the precipitation predicted by the model physical parameterizations in MERRA-2

209     alongside the corrected version in the same fashion as Fig 2. The correction greatly reduces

210     bias, cuts RMSE by one third, slightly improves spatial correlations, and increases the

211    number of stations spanning the diagonal by 28%. CFSR significantly underperforms other

212    reanalyses at FLUXNET2015 locations.

213    Precipitation is among the most difficult quantities for models to simulate. We expect

214    among near surface meteorological variables the lowest correlations and largest coefficient

215    of variation for precipitation.  It also has many observationally-based data sets to choose

216    from, providing a robust estimate of skill to be expected from comparing point

217    measurements to gridded data sets. Figure 2 provides generous thresholds, particularly for

218    correlations, to keep in mind when assessing model simulations of the terms of the surface

219    water and energy balance.  As shown below, correlations of 0.7-0.8 are a challenge for

220    models to attain for precipitation, as well as some other water and energy budget terms.

221    Among near surface meteorology (e.g., temperature and specific humidity) and

222    downward surface fluxes (including shortwave and longwave radiation), precipitation has

223    the greatest small-scale variability on monthly to annual time scales, and is thus the most

224    difficult land surface "forcing" to replicate at the FLUXNET2015 sites. Figures S2-S6 show

225    the scatters and statistics for the models listed in Table 1 for these five variables. Here, the

226    restriction that the years of the models match those at each FLUXNET2015 site is lifted, and

227    the climatologies of the complete data sets are compared. Not surprisingly, the global

228    distribution of annual mean temperature is very well reproduced by the models (Fig. S2),

229    with 88-96% of the observed variance explained. Observed specific humidity is only

230    slightly less well correlated among the models (Fig. S3), but there is a consistent positive

231    bias relative to FLUXNET2015 measurements. Patterns of annual mean downward

232    radiation (Figs. S4 and S5) are well simulated, with a tendency for a slight negative bias in

233    longwave radiation (Fig. S5), and a stronger positive bias in shortwave radiation across

234    models (Fig. S4), consistent with other assessments of model shortwave errors that depend

235    on GCM radiative transfer parameterizations (cf. Slater 2016). Precipitation shows the least

236    agreement; note the bottom row of Fig. S6 is not identical to that of Fig. 2 because the years

237    compared differ. Nevertheless, the results are similar. We can consider MERRA-2 as

238    representing the upper limit of comparison for annual precipitation when the periods do

239    not match between models and observations. Offline Catchment actually performs slightly

240    better than MERRA-2, and CFSv2 is generally the poorest performing model system in the

241    set. Free-running climate models understandably perform worse than either reanalyses or

242    offline LSM simulations, as they are least constrained by observational data. In the case of

243    CFSv2, there are essentially no constraints within the Earth system as an ocean model is

244    coupled; other free-running simulations have specified SSTs.

245       Precipitation is a major source of error at the land surface, but so are elements of the

246    radiation budget. We employ Taylor diagrams to synthesize the statistics of correlation

247    across FLUXNET2015 sites; RMSE and standard deviation are normalized by observed

248    values. Figure 3 shows the global distribution of annual mean downward radiation terms is

249    well simulated across all model configurations, with downward shortwave radiation

250    performing slightly better than downward longwave radiation. Recall for the LSM-only

251    models, downward radiation is an input forcing, and the quality of those data sets can vary

252    significantly (Slater 2016). However, the distribution of upward shortwave radiation is

253    rather poorly simulated, with the NCEP models showing the worst correlations, and the

254    NCAR models the best (yet explaining less than half of the variance). There is also a strong

255    tendency to under-represent the spatial variability (normalized standard deviations less

256    than 1) of downward shortwave radiation. This degrades simulation of net radiation, which

257    has consistently lower correlations than downward radiation terms, yet uniformly better

258    than upward shortwave radiation. The overlap of the spans of annual mean values from

259    models and observations (size of the dots) generally decrease from shortwave down to

260    longwave down to shortwave up.

261       Figure 3 implies discrepancies in the representation of surface albedo across models at

262    FLUXNET2015 sites. We show a Taylor diagram for calculated albedo in Fig. 4. As there are

263    many sites at relatively high northern latitudes that experience snow cover for some part of

264    the year, snow albedo could specifically be a problem. However, a plot of only the JJA

265    albedo verification shows boreal summer generally has even lower fidelity, and

266    systematically low spatial variability, compared to the annual mean. The overlap between

267    the spans of annual mean albedos range among the models from 16% to 38% of

268    FLUXNET2015 sites, but for JJA they span only 13-24%.

269       The low variability could be explained by the fact that most LSMs, whether stand-alone

270    or coupled, have a simple parameterization of albedo based on properties of a small

271    number of vegetation and soil types, often specified as a climatological seasonal cycle. CLM

272    actually calculates surface albedo based on a number of properties including vegetation

273    density and zenith angle of the sun, which may lead to the somewhat better performance of

274    the NCAR models. As described later, the offline NCEP LSM (identified as NL) specifies a

275    multi-year satellite-derived monthly green vegetation fraction as a boundary condition that

276    appears in Fig. 4 to enhance variability, while its positive biases have been noted by Xia et

277    al. (2012). Furthermore, discrepancies between grid box average albedo and local

278    conditions at field sites, including the effect of vegetation differences and soil moisture on

279    albedo (Zaitchik et al. 2013), could add spatial "noise" to the FLUXNET2015 values relative

280    to what models are representing. Nevertheless, such discrepancies lead to a degradation in

281    the representation of surface available energy that is partitioned between sensible, latent

282    and ground heat fluxes. Even an otherwise "perfect" LSM could not produce the right values

283  of these fluxes if net radiation is incorrect. Coupled with errors in precipitation, which

284  affect available soil moisture and thus Bowen ratios, LSMs are at a compounded

285  disadvantage in simulating the surface water and energy budget terms.

286      In Fig. 5 we correlate across the stations the mean errors in key water and energy cycle

287  quantities and present a schematic representation of the relative coupling or

288  connectedness exhibited between terms. This also suggests how errors in the simulation or

289  specification of one term can propagate to others through the land-atmosphere coupling

290  process chain (cf. Santanello et al. 2011). $r_s$ is generally larger than $r_p$ because it does not

291  overemphasize outliers, thus is used for this comparison. Ratios show the fraction of

292  models with correlations at the 90% confidence level, and p-values are based on the

293  average correlation across models. Note the number of included stations varies depending

294  on the availability of observed data (recall from Fig. 1 that a number of FLUXNET2015 sites

295  do not allow for albedo estimations) and among models depending on whether the

296  corresponding grid box is water or land. Furthermore, the data saved from the free-

297  running ECMWF model simulations (EC) do not allow for estimation of albedo, so 11

298  models are compared for albedo.

299      Unsurprisingly, we find surface net radiation errors correlate strongly to albedo errors,

300  with 11 of 11 models registering significant correlations (two-tailed p-values < 0.05) and

301  the multi-model average correlation across 114-118 sites has a p-value of $4\times10^{-7}$. For net

302  radiation versus precipitation, only 2 of 12 models (CL and M1) show significant

303  correlation across 144-151 sites and p=0.55 for the multi-model average, so no direct

304  arrow is drawn in Fig. 5. Note that precipitation errors arise not only from

305  misrepresentation of land-atmosphere interactions, but also from the parameterization of

306  dynamic and thermodynamic processes (so-called "model physics") in the GCM.

12

307    FLUXNET2015 reports both raw and Bowen-ratio corrected heat fluxes. Corrected fluxes

308    are available at fewer than 100 of the sites (two-tailed p=0.05 for correlations $|r| \gtrsim 0.2$,

309    compared to $|r| \gtrsim 0.16$ for the full set of sites), but generally correspond better to the

310    models than uncorrected fluxes, which do not close the surface energy balance (cf. Figs. S9-

311    S12). Regardless, the same story emerges with either set of fluxes: precipitation errors

312    correlate significantly to latent heat flux errors (p=0.02 in Fig. 5) but not sensible heat flux

313    errors (p=0.31). Meanwhile, albedo errors are very strongly linked to sensible heat flux

314    errors (p=7x10$^{-5}$) but not latent heat flux errors (p=0.69). Evaporative fraction (EF; the

315    fraction of sensible + latent heat flux accounted for by the latent heat flux) relates strongly

316    to both, but more strongly to errors in albedo (p=0.003) than precipitation (p=0.05).

317    Consistently, correlating EF errors to the heat flux errors (black two-way arrows)

318    demonstrates more variance explained by sensible heat flux than latent heat flux. Finally,

319    LCL errors relate strongly to precipitation errors (p=2x10$^{-5}$) but are marginally significant

320    in relation to albedo errors (p=0.06). LCL has a prevalent negative bias (Fig. S8) reflecting

321    the positive biases in specific humidity.

322    This analysis shows that models have troublesome errors in both the surface water and

323    energy cycles, which make their way into the land-atmosphere coupling process chain. As a

324    result, the degree to which weather and climate models correctly simulate feedbacks of

325    land surface anomalies onto the atmosphere may be cast into some doubt. However, the

326    origins of several sources of error have been identified and their alleviation can be

327    pursued. In section 5 we will examine directly model fidelity in simulating metrics of land-

328    atmosphere coupling.

329

330    **4. Mean Annual Cycle**

331     The next criterion for models, beyond simulating the annual means among

332     FLUXNET2015 sites, is reproducing the annual cycle. The first harmonic is fit to the 12

333     monthly means for each variable, determining phase and magnitude (half of valley-to-peak

334     distance) using a standard Fourier transform. Errors in phase and magnitude at each

335     station, quantified across all stations with similar metrics as the annual mean, indicate skill

336     in simulating the annual cycle. Amplitude errors are displayed in conventional scatter

337     diagrams (see Figs. S15-S24), but to display information for phase errors, we have

338     configured the classical scatter diagram in a polar projection (see Figs. S25-34; the caption

339     of Fig. S25 gives a detailed description of those plots). The whiskers in the supplemental

340     figures again show models frequently display a smaller range of year-to-year variability

341     than data from FLUXNET2015 sites. This may be partially explained by the scale difference

342     (point measurements will vary more than grid-box averages) but is also likely due to the

343     overly deterministic nature of many model parameterizations (Palmer 2012).

344     Taylor diagrams summarize the results across models. We focus on depictions of energy

345     budget terms, as they reveal some of the main issues among models. Figure 6 shows model

346     performance in simulating the amplitudes of the annual cycles of net radiation, sensible

347     and latent heat fluxes across FLUXNET2015 sites. All model products demonstrate similar

348     skill for net radiation, clustered between 0.64-0.78 correlation and a tendency toward too

349     large an annual cycle. Only the offline NCEP and coupled ECMWF models have a negative

350     bias in amplitude. Latent heat flux simulations show lower skill for every model, clustering

351     between 0.28-0.43 for correlations. At the stations where energy balance corrected fluxes

352     are provided, correlations improve to 0.37-0.50 (not shown). The positive bias is not so

353     pervasive for latent heat; rather it appears the positive bias in net radiation tends to be

354     expressed in the sensible heat term. There is also a much larger spread among models for

355  sensible heat, both in terms of correlation (0.14-0.54) and normalized standard deviation

356  (0.78-1.50).

357      The models' skill in representing the phase of the annual cycle has a similar distribution

358  (Fig. 7). The phase of net radiation is best represented, latent and sensible heat have spatial

359  correlations of phasing between ~0.8-0.92 with sensible heat phases having slightly lower

360  fidelity in general. It is interesting as the general consensus is that sensible heat flux is a

361  simpler process to model than latent heat flux, yet it has been shown in other contexts that

362  LSMs struggle more to simulate sensible heat flux (e.g., Best et al. 2015).

363      The Taylor diagram for the annual cycle of albedo (Fig. 8) shows very similar

364  correlations of the yearly amplitude between models and observations (0.50-0.71) but a

365  large range in standard deviation;  Noah v2.7.1 (NL) shows a particularly high value

366  contributing to large RMSE. The phase is better represented by all models, but interestingly

367  the standard deviations are uniformly over-estimated. Most models now use global MODIS-

368  based data sets of albedo as either a parameter set or for calibration of surface radiative

369  parameterizations, so the large inter-model spread and lack of obvious clustering within

370  families of models is surprising.

371

372  **5. Coupling Metrics**

373      Correlations between land surface state variables and surface fluxes (the terrestrial leg

374  of coupling) and between land surface fluxes and atmospheric states or properties

375  (atmospheric leg) may indicate feedbacks. For instance in the terrestrial leg, positive

376  (negative) correlation between soil moisture and latent (sensible) heat flux implies soil

377  moisture control of fluxes (a moisture limited situation) as opposed to energy (net

378  radiation) limited situations where atmospheric states control the fluxes. However, the

15

379    variance in the driving term(s) must also be sufficiently large for a sensitivity of

380    atmosphere to the land to have a consequential impact on climate, relative to other factors.

381    A coupling index $I$ can be constructed from terms in either leg: $I = \sigma(b)r(a,b) = \sigma(a)\frac{db}{da}$

382    where $a$ is the forcing and $b$ is the responding variable, $\sigma$ is standard deviation in time, $r$ is

383    correlation in time, and the linear regression slope of $b$ on $a$ is a measure of the sensitivity

384    of $b$ to $a$ (Dirmeyer 2011, Dirmeyer et al. 2013).

385    Figure 9 synthesizes the performance of the various model configurations regarding

386    two-legged coupling metrics linking soil moisture to boundary layer properties. The

387    formulae for the coupling indices are indicated on the figure axes calculated from daily

388    mean values. The terrestrial leg quantifies the combined sensitivity (correlation) of surface

389    fluxes (here, latent heat flux) to land states (soil moisture) with variability (standard

390    deviation) of the flux. The atmospheric leg links surface fluxes (sensible heat flux) to

391    atmospheric states (LCL, which combines near surface temperature and humidity

392    information). Larger values denote stronger feedback linkages.

393    In each panel of Fig. 9, similar to the approach of Sippel et al. (2017), quantities are

394    calculated for the three consecutive months that have the warmest average temperature

395    according to the FLUXNET2015 data. We distinguish between positive values of each

396    metric, which indicate the existence of feedbacks from land to atmosphere, from negative

397    (no feedbacks) by coloring the four quadrants by their coupling regimes: red = both legs

398    present and a full coupling pathway; green = the land leg is present, the atmospheric leg is

399    missing; blue = atmospheric leg is present, land is missing; grey = neither leg present. The

400    white dots show where FLUXNET2015 sites fall in this two-dimensional metric space. The

401    colored dots are each model's rendering of the metrics for the grid boxes containing the

402    FLUXNET2015 sites; the color indicates the quadrant according to the FLUXNET

16

403     measurements. Thus, the more colored dots that fall in the quadrant with the matching

404     color, the better the model is reproducing the global pattern of coupling regimes.

405       The model centroid usually lies below and to the right of the observed centroid for a

406     given coupling regime, meaning models tend to over-estimate the terrestrial coupling index

407     (the rightward offset), yet underestimate the strength of the atmospheric leg (the

408     downward offset). Recall the number of FLUXNET2015 sites compared is not the same for

409     each model. The percentage in each quadrant indicates how many of the FLUXNET2015

410     sites in that regime are correctly placed in the right quadrant. For instance, the CFS

411     Reanalysis has 76% of the FLUXNET stations exhibiting both coupling legs (red) in the

412     correct regime. However, there are clearly many dots of other colors also in the red

413     quadrant, showing the model places many other stations erroneously in that regime.

414     Interestingly, none of the models put the few sites with no warm-season coupling in the

415     grey quadrant. Overall, the reanalyses perform best: a 56.5% overall hit rate for the fully-

416     coupled regime versus 52.8 for coupled models, and 44.0% for offline LSMs; and for the

417     atmosphere-only coupling regime 49.2% versus 33.0% for coupled models and 31.6% for

418     offline LSMs.

419       We have also examined performance of the models for their simulation of the observed

420     FLUXNET2015 correlations and standard deviations (the two terms in the coupling

421     indices) separately. As implied previously for the terrestrial leg, there is a positive bias in

422     correlations for all models except for ERA-Interim (Table 2). Bias in the standard deviation

423     of latent heat fluxes across all sites is small for most models, so most of the positive bias in

424     coupling index comes from the correlation term. The model biases are even stronger in the

425     anti-correlation between soil moisture and sensible heat flux (not shown). However, there

426     is generally an even greater bias in correlations for the atmospheric leg (Table 2) paired in

427    every model with an underrepresentation of the daily variability of the LCL. These two

428    biases compound, leading to the strong underrepresentation of coupling in the atmospheric

429    leg of land-atmosphere interactions.

430        There are several caveats to note. First, the notion of calculating the atmospheric

431    coupling leg from offline LSM simulations is only partially justifiable. It is certainly possible

432    to calculate the correlations between surface fluxes and LCL height (which depends on

433    near-surface meteorological data supplied as forcing to the LSM), but there is no possibility

434    for the fluxes to affect 2m temperature or humidity. Thus, this is more of a test of model

435    consistency than a true diagnosis of coupling.

436        Second, estimates of the correlation component of the coupling indices from observed

437    data must be closer to zero than the true values in nature, because random measurement

438    errors will degrade correlations (Robock et al. 1995). Thus, it is not necessarily wrong that

439    models show a stronger terrestrial coupling leg than FLUXNET2015 data. The degree of

440    impact can be estimated for variables such as soil moisture, whose auto-correlation time

441    scales are much longer than the daily data interval (cf. Dirmeyer et al. 2016) but can be

442    difficult to estimate from small samples or for other quantities. Nevertheless, the fact that

443    models routinely underestimate the strength of the atmospheric leg runs counter to being

444    attributable to random observational errors at FLUXNET sites, and likely represents real

445    model bias.

446        Finally, the difference in scale between flux tower measurements (typically

447    representative of conditions in an area of a square kilometer or less) and model grid-box

448    averages (here ranging from $200–2x10^4$ km$^{-2}$) can affect statistics. Dirmeyer et al. (2016)

449    showed there was little sensitivity of estimates of temporal variations in daily soil moisture

450    to spatial scale differences in the model grid box range, however, the same may not be true

451     for other terms, or for correlations. The larger the averaging area, the smoother we should

452     expect time series to be, potentially affecting estimation of coupling indices.

453

454     **6. Discussion and Summary**

455         We have confronted four different global model systems in multiple configurations (LSM

456     only, LSM coupled to GCM, and reanalysis) with flux tower observations from 166 sites in

457     the global FLUXNET2015 data set to determine how well they reproduce the spatial

458     distribution of annual means and the annual cycle of state variables and terrestrial surface

459     fluxes, and coupling indices between land and atmosphere. Returning to Table 2, there is a

460     separation evident between the three classes of models. For the terrestrial leg of land-

461     atmosphere coupling, all models appear to overestimate correlations between soil

462     moisture and latent heat flux, with the caveat discussed previously that correlations

463     necessarily skew low when calculated from observed data. Nevertheless, assuming as much

464     as a 50% reduction from true correlations, it appears the reanalyses do the best job at

465     reproducing observed correlations, followed by the free-running models and last the

466     uncoupled LSMs. There is a similar stratification for the standard deviation of latent heat

467     flux: reanalyses very closely represent the observed temporal variability of this flux, while

468     coupled models and stand-alone LSMs progressively underestimate it. For the atmospheric

469     leg, represented by the coupling index between sensible heat flux and LCL height, all

470     classes of models severely underestimate the correlation and the day-to-day variability in

471     the LCL. Reanalyses again do the best job at correlations and stand-alone LSMs are the

472     worst. Here, coupled models fare slightly better than reanalyses in representing LCL

473     variance. Given that reanalyses are somewhat constrained by the assimilation of

474     observations, the errors in those models do not manifest as freely, so it makes sense

475  reanalyses should verify the best. On the other hand, offline LSMs lack some of the coupling

476  we are trying to gauge. For example, surface sensible and latent heat fluxes cannot affect

477  near surface temperature and humidity in such a configuration. This prescription of near-

478  surface states interferes with the feedback processes.

479  General characteristics of note are that scatter diagrams of model versus FLUXNET2015

480  quantities almost always show a linear regression slope indicating a wider range of

481  variation in the observations. Models also tend to have lower interannual variability

482  (length of whiskers) than observations suggest. These traits are consistent with scale

483  differences between model grid cells and the area sampled by flux towers; model grid

484  values represent areas at least 2-4 orders of magnitude larger, which particularly affects

485  precipitation forcing. Thus, this difference is not a concern regarding model performance

486  *per se*, but rather representativeness across scales.

487  Another general characteristic is that the models verify better against the corrected

488  surface fluxes and quantities derived from them; wherein observed sensible and latent heat

489  values are adjusted to close the surface energy budget. This makes sense as models close

490  surface energy (and water) budgets by design, whereas closure is not assured in an

491  observational setting where a number of instruments, with different calibrations and error

492  characteristics, contribute separate terms of the surface balances. However, when the

493  propagation of model errors through the energy and water cycles are traced (Fig. 5), EF in

494  models shows strong sensitivity to radiation errors, implying that conservation of Bowen

495  ratio (and thus EF) as a means to correct observed heat fluxes and close the energy balance

496  may not be the most efficacious.

497  There are differences that do appear to reflect general model biases. All models and

498  configurations show a positive bias in near-surface humidity (Fig. S3, S14), downward

499    shortwave radiation (Figs. S4, S17) and a range of biases in downward longwave radiation

500    (Fig. S5). Such radiation biases are a long-standing problem in global models (cf. Dirmeyer

501    et al. 2006), and stem from problems in the parameterization of atmospheric radiative

502    transfer, clouds and aerosols in GCMs. However, not all radiative errors are atmospheric in

503    origin – there is clear indication that LSMs struggle to represent the spatial and temporal

504    variability of surface albedo (Figs. 4, 8).

505    Combined with well-known difficulties models have in simulating precipitation (Figs. 2,

506    S6, S15, S25), it becomes extremely challenging for models to partition available energy

507    correctly at the surface between latent, sensible and ground heat fluxes, and to reproduce

508    the spatiotemporal patterns of relationships between soil moisture, surface fluxes and the

509    lower troposphere. Errors in latent heat flux generally correlate significantly to

510    precipitation errors, while sensible heat flux errors relate strongly to surface albedo errors.

511    Evaporative fraction errors connect to both, but more strongly to the energy (albedo –

512    sensible heat flux) pathway than the water (precipitation – latent heat flux) pathway.

513    Height of the LCL, which has a strong negative bias across all models related to the positive

514    humidity bias, has errors that correlate strongly to the water cycle pathway, but also to the

515    energy cycle pathway.

516    The spatial distributions of the annual cycles are generally well reproduced for energy

517    budget terms, except for upward shortwave radiation, related to the albedo problems

518    discussed earlier. However, there is a tendency for too strong a seasonal cycle in net

519    radiation, caused by excessive summertime downward shortwave radiation, and expressed

520    more strongly in the annual cycle of sensible heat flux than latent heat flux. Models

521    generally do very well representing the spatial distribution of the phasing of the annual

21

522  cycle, even for precipitation (64-92% of variance explained) and soil moisture (40-61% of

523  variance explained).

524      Finally, despite the barriers described above to models' capacity to represent the

525  spatiotemporal distribution of land-atmosphere coupling, we find models often do a

526  reasonable job. Some systematic biases are evident: models consistently over-estimate the

527  strength of the terrestrial leg of coupling (namely, too strong a correlation between soil

528  moisture and sensible heat fluxes), yet even more clearly underestimate the strength of the

529  atmospheric leg (both the correlation between surface fluxes and boundary layer

530  properties, and day-to-day variability of boundary layer properties). Random

531  observational error tends to reduce correlations between observed quantities, so it is

532  possible that models are not greatly overestimating the terrestrial leg of coupling, or

533  perhaps are not overestimating it at all. However, we find the time series at most

534  FLUXNET2015 sites are too short to robustly estimate the random error effects on

535  correlation – perhaps in another ten years we will be able to quantify these errors.

536  Similarly, the spatial scale differences between observations and model output may

537  contribute to the variance differences in the atmospheric leg, but disparity in correlations

538  between surface fluxes and LCL could only be stronger than calculated here, not weaker,

539  because of the effect of measurement error.

540      LSMs forced by global gridded meteorology rather than local forcing from the tower

541  sites themselves are handicapped to some degree (cf. Chen et al. 2017). So our most

542  confident conclusion regarding land-atmosphere coupling is that models under-represent

543  the feedback of surface fluxes on boundary layer properties at FLUXNET2015 sites. We find

544  this unique data set has potential for model development and parameter optimization to

545      alleviate biases in model configurations shown to mirror those used in forecasting

546      applications (Orth et al. 2016, 2017).

547      Overall, we conclude that many of the long-known problems and biases in global models

548      of the land-atmosphere portion of the climate system still exist. Nevertheless, there is a fair

549      degree of compensation among errors, such that model representations of land-

550      atmosphere coupling often appear fairly good. Some targets for model improvement are

551      clear, however, as coupling linkages suggest processes where problems may lie. The

552      representation of surface albedo (LSM) and the quantities of downward radiation at the

553      surface (GCM) need improvement among the energy cycle terms, along with the

554      partitioning of available energy between latent and sensible heat flux (a coupled model

555      development problem). Precipitation errors remain large, and inconsistencies in

556      representing soil moisture among models and between models and nature (cf. Koster et al.

557      2009) remain stubborn issues.

558      As one might expect, reanalyses tend to perform better, as they are more constrained by

559      observational data. LSMs run offline also benefit from meteorological forcing that is highly

560      observational in origin, but can be handicapped by their lack of two-way interaction with

561      the lower troposphere. It should be clear from the various figures that individual models

562      perform better or worse at simulating specific facets of land-atmosphere interactions.

563      However, we emphasize here the commonalities among models more than differences. This

564      study is not primarily intended as a model inter-comparison, but rather a multi-model

565      attempt to draw model-independent conclusions about the current state of performance of

566      land-atmosphere models (in various configurations) by confronting them with a new and

567      unique observational data set.

568    Furthermore, this study is not a final judgement, but a first look that will hopefully

569    catalyze accelerated development and improvement in coupled land-atmosphere modeling.

570    Application of cross-component metrics like coupling indices can reveal prime areas for

571    model development that are not evident from piecewise evaluation of model components.

572    The next step is intensive, focused sensitivity studies with individual models, preferably

573    validated in the context of coupled model systems, that will zero in on the problematic

574    parameterizations. We may also need to revisit some of the fundamental assumptions that

575    underpin the formulations in models (e.g., Cheng et al. 2017).

576    Furthermore, it is clear that long-term observational monitoring is highly valuable, and

577    that value only increases with the duration of data sets at individual sites. Greater spatial

578    distribution of flux tower sites, especially into under-monitored regions outside middle-

579    and high-latitudes, would further increase the overall usefulness to model development.

580

**References:**

Andreae, M. O., and co-authors, 2002: Biogeochemical cycling of carbon, water, energy, trace gases, and aerosols in Amazonia: The LBA-EUSTACH experiments. *J. Geophys. Res.*, **107**, LBA33-1–LBA33-25.

Baldocchi, D., and co-authors, 2001: FLUXNET: A new tool to study the temporal and spatial variability of ecosystem-scale carbon dioxide, water vapor and energy flux densities. *Bull. Amer. Meteor. Soc.*, **82**, 2415-2434.

Balsamo, G., and co-authors, 2015: ERA--Interim/Land: a global land surface reanalysis data set, *Hydrol. Earth Syst. Sci.*, **19**, 389- 407, doi: 10.5194/hess-19-389- 2015.

Balzarolo, M., and co-authors, 2014: Evaluating the potential of large-scale simulations to predict carbon fluxes of terrestrial ecosystems over a European Eddy Covariance network, *Biogeosci.*, **11**, 2661-2678, doi:10.5194/bg-11-2661-2014.Best, M. J., and Co-authors, 2015: The plumbing of land surface models: benchmarking model performance. *J. Hydrometeor.*, **16**, 1425-1442, doi: 10.1175/JHM-D-14-0158.1.

Bonan, G. B., K. W. Oleson, R. A. Fisher, G. Lasslop, and M. Reichstein, 2012: Reconciling leaf physiological traits and canopy flux data: Use of the TRY and FLUXNET databases in the Community Land Model version 4, *J. Geophys. Res.*, **117**, G02026, doi: 10.1029/2011JG001913.

Boussetta, S., G. Balsamo, A. Beljaars, T. Kral, L. Jarlan, 2013: Impact of a satellite-derived leaf area index monthly climatology in a global numerical weather prediction model. *Int. J. Remote Sens.*, **34**, 3520-3542. doi: 10.1080/01431161.2012.716543.

Chen, L., P. A. Dirmeyer, Z. Guo and N. M. Schultz, 2017: Pairing FLUXNET sites to validate model representations of land use/land cover change. *Hydrol. Earth Sys. Sci. Discus.*, doi: 10.5194/hess-2017-190.

Cheng, Y., C. Sayde, Q. Li, J. Basara, J. Selker, E. Tanner, and P. Gentine, 2017: Failure of Taylor's hypothesis in the atmospheric surface layer and its correction for eddy-covariance measurements. *Geophys. Res. Lett.*, **44**, 4287–4295, doi: 10.1002/2017GL073499.

624 Dee, D. P., and co-authors, 2011: The ERA-Interim reanalysis: configuration and
625    performance of the data assimilation system. *Quart. J. Roy. Meteor. Soc.*, **137**, 553-597,
626    doi: 10.1002/qj.828.

627 Dirmeyer, P. A., R. D. Koster, and Z. Guo, 2006: Do global models properly represent the
628    feedback between land and atmosphere? *J. Hydrometeor.*, **7**, 1177-1198, doi:
629    10.1175/JHM532.1.

630 Dirmeyer, P. A., 2011: The terrestrial segment of soil moisture-climate coupling. Geophys.
631    Res. Lett., 38, L16702, doi: 10.1029/2011GL048268.

632 Dirmeyer, P. A., S. Kumar, M. J. Fennessy, E. L. Altshuler, T. DelSole, Z. Guo, B. Cash and D.
633    Straus, 2013: Model estimates of land-driven predictability in a changing climate from
634    CCSM4. J. Climate, 26, 8495-8512, doi: 10.1175/JCLI-D-13-00029.1.

635 Dirmeyer, P. A., and co-authors, 2016: Confronting weather and climate models with
636    observational data from soil moisture networks over the United States. *J. Hydrometeor.*,
637    **17**, 1049-1067, doi: 10.1175/JHM-D-15-0196.1.

638 Dirmeyer, P. A., P. Gentine, M. B. Ek, and G. Balsamo, 2017: Land Surface Processes Relevant
639    to S2S Prediction. [Chapter 8 in: *The Gap Between Weather and Climate Forecasting:*
640    *Sub-Seasonal to Seasonal Prediction* (A. W. Robertson and F. Vitart Eds.)], Elsevier, (in
641    revision).

642 Dorigo, W. A., and co-authors, 2011: The International Soil Moisture Network: a data
643    hosting facility for global in situ soil moisture measurements, *Hydrol. Earth Syst. Sci.,* **15**,
644    1675-1698, doi: 10.5194/hess-15-1675-2011.

645 Dorigo, W.A., and co-authors, 2013: Global automated quality control of in situ soil
646    moisture data from the International Soil Moisture Network. *Vadose Zone J.*, **12**(3), doi:
647    10.2136/vzj2012.0097.

648 Dorigo, W., and co-authors, 2017: ESA CCI Soil Moisture for improved Earth system
649    understanding: state-of-the art and future directions, *Remote Sens. Env.* (in press),
650    10.1016/j.rse.2017.07.001.

651     Ek, M. B., K. E. Mitchell, Y. Lin, E. Rogers, P. Grunmann, V. Koren, G. Gayno, and J. D. Tarplay,
652        2003: Imple- mentation of Noah land surface model advances in the National Centers for
653        Environmental Prediction operational mesoscale Eta model. *J. Geophys. Res.*, **108**, 8851,
654        doi: 10.1029/2002JD003296.

655     Famiglietti, J. S., J. A. Devereaux, C. A. Laymon, T. Tsegaye, P. R. Houser, T. J. Jackson, S. T.
656        Graham, M. Rodell, and P. J. van Oevelen, 1999: Ground-based investigation of soil
657        moisture variability within remote sensing footprints during the Southern Great Plains
658        97 (SGP97) hydrology experiment. *Water Resour. Res.*, **35**, 1839-1851.

659     Gelaro, R., and co-authors, 2017: The Modern-Era Retrospective analysis for Research and
660        Applications, version 2 (MERRA-2). *J. Climate*, **30**, 5419-5454, doi: 10.1175/JCLI-D-16-
661        0758.1.

662     Global Modeling and Assimilation Office (GMAO), 2015: MERRA-2 inst1_2d_lfo_Nx: 2d, 1-
663        Hourly, Instantaneous, Single-Level, Assimilation, Land Surface Forcings V5.12.4,
664        Greenbelt, MD, USA, Goddard Earth Sciences Data and Information Services Center (GES
665        DISC), Accessed 3 July 2016, doi: 10.5067/RCMZA6TL70BG.

666     Global Modeling and Assimilation Office (GMAO), 2015: MERRA-2 tavg1_2d_lfo_Nx: 2d, 1-
667        Hourly, Time-Averaged, Single-Level, Assimilation ,Land Surface Forcings V5.12.4,
668        Greenbelt, MD, USA, Goddard Earth Sciences Data and Information Services Center (GES
669        DISC), Accessed 3 July 2016, doi: 10.5067/L0T5GEG1NYFA.

670     Jackson, T. J., and A. Y. Hsu, 2001: Soil moisture and TRMM microwave imager relationships
671        in the Southern Great Plains 1999 (SGP99) experiment, IEEE Trans. *Geosci. Remote
672        Sens.*, **39**, 1632-1642.

673     Kinter III, J. L., and co-authors, 2013: Revolutionizing climate modeling – Project Athena: A
674        multi-institutional, international collaboration. *Bull. Amer. Meteor. Soc.*, **94**, 231–245,
675        doi: 10.1175/BAMS-D-11-00043.1.

676     Koster, R. D., Z. Guo, P. A. Dirmeyer, R. Yang, K. Mitchell, and M. J. Puma, 2009: On the nature
677        of soil moisture in land surface models. *J. Climate*, **22**, 4322–4335, doi:
678        10.1175/2009JCLI2832.1.

679    Lawrence, D. M., and co-authors, 2011: Parameterization improvements and functional and

680        structural advances in version 4 of the Community Land Model. *J. Adv. Model. Earth*

681        *Syst.*, **3**, doi: 10.1029/2011MS000045.

682    Mahanama, S. P. P., R. D. Koster, G. K. Walker, L. L. Takacs, R. H. Reichle, G. De Lannoy, Q. Liu,

683        B. Zhao, and M. J. Suarez, 2015: Land Boundary Conditions for the Goddard Earth

684        Observing System Model Version 5 (GEOS-5) Climate Modeling System - Recent Updates

685        and Data File Descriptions. NASA/TM–2015-104606, Vol. 39, 55 pp. Document (4608

686        kB).

687    Melaas, E. K., A. D. Richardson, M. A. Friedl, D. Dragoni, C. M. Gough, M. Herbst, L.

688        Montagnani, and E. Moors, 2013: Using FLUXNET data to improve models of springtime

689        vegetation activity onset in forest ecosystems. *Ag. Forest Meteor.*, **171-172**, 46-56.

690    Mitchell, K., 2005: The Community Noah Land Surface Model User's Guide Public Release

691        Version                    2.7.1,                    [available                    at:

692        http://www.ral.ucar.edu/research/land/technology/lsm/noah/Noah_LSM_USERGUIDE

693        _2.7.1.pdf].

694    Molod, A., Takacs, L., Suarez, M., and Bacmeister, J., 2015: Development of the GEOS-5

695        atmospheric general circulation model: evolution from MERRA to MERRA2, *Geosci.*

696        *Model Dev.*, **8**, 1339-1356, doi: 10.5194/gmd-8-1339-2015.

697    Orth, R., E. Dutra, and F. Pappenberger, 2016: Improving weather predictability by

698        including land surface model parameter uncertainty. *Mon. Wea. Rev.,* **144**, 1551–1569,

699        doi: 10.1175/MWR-D-15-0283.1.

700    Orth, R., Dutra, E., Trigo, I. F., and Balsamo, G., 2017: Advancing land surface model

701        development with satellite-based Earth observations, *Hydrol. Earth Syst. Sci.*, **21**, 2483-

702        2495, doi:10.5194/hess-21-2483-2017.

703    Palmer, T. N., 2012: Towards the probabilistic Earth-system simulator: a vision for the

704        future of climate and weather prediction. *Quart. J. Roy. Meteor. Soc.*, **138**, 841-861.

705 Pastorello, G. Z., D. Papale, H. Chu, C. Trotta, D. A. Agarwal, E. Canfora, D. D. Baldocchi, and
706     M. S. Torn, 2017: A new data set monitors land-air exchanges. *EOS Earth & Space Science*
707     *News*, **98**(8), 28-32.

708 Peters-Lidard, C. D., and co-authors, 2007: High performance earth system modeling with
709     NASA/GSFC's Land Information System. *Innov. Syst. Software Eng.*, **3**, doi:
710     10.1007/s11334-007-0028-x.

711 Purdy, A. J., J. B. Fisher, M. L. Goulden, and J. S. Famiglietti, 2016. Ground heat flux: An
712     analytical review of 6 models evaluated at 88 sites and globally. *J. Geophys. Res.,* **121**,
713     3045-3059.

714 Quiring, S. M., T. W. Ford, J. K. Wang, A. Khong, E. Harris, T. Lindgren, D. W. Goldberg, and Z.
715     Li, 2016: North American Soil Moisture Database: Development and applications. *Bull.*
716     *Amer. Meteor. Soc.*, **97**, 1441-1460.

717 Reichle, R. H., and Q. Liu, 2014. Observation-Corrected Precipitation Estimates in GEOS-5.
718     NASA/TM–2014-104606, Vol. 35. http://gmao.gsfc.nasa.gov/pubs/docs/Reichle734.pdf.

719 Reichle, R. H., C. S. Draper, Q. Liu, M. Girotto, S. P. Mahanama, R. D. Koster, and G. De Lannoy,
720     2017a. Assessment of MERRA-2 land surface hydrology estimates. *J. Climate*, **30**, 2937-
721     2960, doi: 10.1175/JCLI-D-16-0720.1.

722 Reichle, R., Q. Liu, R. Koster, C. Draper, S. Mahanama, and G. Partyka, 2017b. Land surface
723     precipitation in MERRA-2. *J. Climate*, **30**, 1643-1664, doi: 10.1175/JCLI-D-16-0570.1.

724 Reichstein, M., and co-authors, 2005: On the separation of net ecosystem exchange into
725     assimilation and ecosystem respiration: review and improved algorithm. *Glob. Change*
726     *Biol.*, **11**, 1424-1439, doi: 10.1111/j.1365-2486.2005.001002.x.

727 Rienecker, M. M., and co-authors, 2011: MERRA: NASA's Modern-Era Retrospective
728     Analysis for Research and Applications. *J. Climate*, **24**, 3624-3648, doi:10.1175/JCLI-D-
729     11-00015.1.

730 Robock, A., K. Ya. Vinnikov, C. A. Schlosser, N. A. Speranskaya and Y. Xue, 1995: Use of
731     midlatitude soil moisture and meteorological observations to validate soil moisture
732     simulations with biosphere and bucket models.. *J. Climate*, **8**, 15-35.

733  Saha, S., and co-authors, 2010: The NCEP Climate Forecast System Reanalysis. *Bull. Amer.*
734      *Meteor. Soc.*, **91**, 1015–1057, doi: 10.1175/2010BAMS3001.1.

735  Santanello, J. A., C. D. Peters-Lidard, and S. V. Kumar, 2011: Diagnosing the sensitivity of
736      local land-atmosphere coupling via the soil moisture-boundary layer interaction. *J.*
737      *Hydrometeor.*, **12**, 766-786.

738  Santanello, J. A., P. A. Dirmeyer, C. R. Ferguson, K. L. Findell, A. B. Tawfik, A. Berg, M. B. Ek, P.
739      Gentine, B. Guillod, C. van Heerwaarden, J. Roundy, and V. Wulfmeyer, 2017: Land-
740      atmosphere interactions: The LoCo perspective. *Bull. Amer. Meteor. Soc.*, (in revision).

741  Sellers, P. J., F. G. Hall, G. Asrar, D. E. Strebel, and R. E. Murphy, 1992: An overview of the
742      First International Satellite Land Surface Climatology Project (ISLSCP) Field Experiment
743      (FIFE). *J. Geophys. Res.*, **97**, 18,345-18,372.

744  Sellers, P. J., and co-authors, 1995: The Boreal Ecosystem-Atmosphere Study (BOREAS): An
745      overview and early results from the 1994 field year. *Bull. Amer. Meteor. Soc.*, **76**, 1549-
746      1577.

747  Sheffield, J., G. Goteti, and E. F. Wood, 2006: Development of a 50-yr high-resolution global
748      dataset of meteorological forcings for land surface modeling. *J. Climate*, **19**, 3088-3111.

749  Shukla, R. P., B. Huang, L. Marx, J. L. Kinter and C.-S. Shin, 2017: Predictability and
750      prediction of Indian summer monsoon by CFSv2: implication of the initial shock effect.
751      *Climate Dy*n. (published online), doi: 10.1007/s00382-017-3594-0.

752  Sippel, S., J. Zscheischler, M. D. Mahecha, R. Orth, M. Reichstein, M. Vogel, and S. I.
753      Seneviratne, 2017: Refining multi-model projections of temperature extremes by
754      evaluation against land–atmosphere coupling diagnostics. *Earth Sys. Dyn.*, **8**, 387-403.

755  Slater, A. G., 2016: Surface solar radiation in North America: A comparison of observations,
756      reanalyses, satellite, and derived products. *J. Hydrometeor.*, **17**, 401-420.

757  Viovy, N., 2013. CRUNCEP data set for 1901–2010, [Available at
758      https://www.earthsystemgrid.org/dataset/ucar.cgd.ccsm4.CRUNCEP.v4.html].

759     Vuichard, N., and D. Papale, 2015: Filling the gaps in meteorological continuous data
760         measured at FLUXNET sites with ERA-Interim reanalysis. *Earth Sys. Sci. Data*, **7**, 157-
761         171, doi: 10.5194/essd-7-157-2015.

762     Williams, M., and co-authors, 2009: Improving land surface models with FLUXNET data.
763         *Biogeosci.* **6**, 1341-1359.

764     Xia, Y., and co-authors, 2012: Continental-scale water and energy flux analysis and
765         validation for the North American Land Data Assimilation System project phase 2
766         (NLDAS-2): 1. Intercomparison and application of model products, *J. Geophys. Res.*, **117**,
767         D03109, doi:10.1029/2011JD016048.

768     Xie, P., and P. A. Arkin, 1997: Global precipitation: A 17-year monthly analysis based on
769         gauge observations, satellite estimates, and numerical model outputs. *Bull. Amer.*
770         *Meteor. Soc.*, **78**, 2539-2558.

771     Xie, P., M. Chen, A. Yatagai, T. Hayasaka, Y. Fukushima, and S. Yang, 2007: "A gauge-based
772         analysis of daily precipitation over East Asia." *J. Hydrometeor.*, **8**, 607–626.

773     Zaitchik, B., F., J. A. Santanello, S. V. Kumar, and C. D. Peters-Lidard, 2013: Representation of
774         soil moisture feedbacks during drought in NASA Unified WRF (NU-WRF). *J.*
775         *Hydrometeor.*, **14**, 360-367.

776 Table 1. Specifications for the four land and atmosphere model systems, including time
777 span of data and spatial resolution. Two-letter abbreviations are used in subsequent
778 figures and tables; generally for the first letter: N=NCEP, M=NASA (MERRA system),
779 C=NCAR (Community models), E=ECMWF; for the second letter: L=LSM run "offline",
780 C=LSM coupled to GCM, R=reanalysis (except that two MERRA reanalyses are included, so
781 they are labeled 1 and 2).

| System | Offline LSM | Free-Running | Reanalysis |
|---|---|---|---|
| **NOAA/ NCEP** | **NL:** Noah2.7.1 [1982-2010] 1°x1° with forcing from Sheffield et al. (2006) | **NC:** CFSv2 [48 years] ~0.94°x0.94° fully coupled Shukla et al. (2017) | **NR:** CFSR [1979-2009] 0.31°x0.37° Saha et al. (2010) |
| **NASA/ GMAO** | **ML:** Catchment with boundary conditions from Mahanama et al (2015) plus physics changes [1980-2015] 0.625°x0.5° with MERRA-2 forcing and corrected precipitation Reichle et al. (2017b), GMAO (2015a,b) | **MC:** GEOS5 Heracles-5 4 p3-M3; LSM as in **ML** [2000-2015] 0.5°x0.5° with observed SST | **M2:** MERRA-2 [1980-2015] 0.625°x0.5° Gelaro et al. (2017); **M1:** MERRA [1980-2015] 0.667°x0.5° Rienecker et al (2011) |
| **NCAR** | **CL:** CLM4.5 [1991-2010] 1.25°x0.9° with CRUNCEP (Viovy 2013) forcing Lawrence et al. (2011) | **CC:** CESM 1.2.2 (CAM4 + CLM4.5) [1991-2014] 1.25°x0.9° with climatological SST | --none-- |
| **ECMWF** | **EL:** HTESSEL 43R1 [1979-2015] TCo639 16km Balsamo et al. (2015) | **EC:** IFS in Athena Project [1961-2007] T1279 interpolated to N80 1.125°x1.125° with observed SST Kinter et al. (2013) | **ER:** ERA-Interim [1979-2015] 0.75°x0.75° Dee et al. (2011) |

782
783

784 Table 2: The average value of the two terms used to calculate the terrestrial and
785 atmospheric coupling indices using data from FLUXNET2015, each model, and averages
786 from various groupings of the models.
787

| | Terrestrial | | Atmospheric | |
|---|---|---|---|---|
| | r(SM,LHF) | $\sigma$(LHF) | r(SHF,LCL) | $\sigma$(LCL) |
| FLUXNET2015 | **0.07** | **21.2** Wm$^{-2}$ | **0.35** | **432** m |
| NL | 0.31 | 18.2 | -0.22 | 221 |
| NC | 0.21 | 21.5 | 0.13 | 412 |
| NR | 0.22 | 23.1 | 0.21 | 396 |
| ML | 0.14 | 15.9 | 0.08 | 366 |
| MC | 0.13 | 14.0 | 0.02 | 291 |
| M2 | 0.11 | 21.4 | 0.12 | 287 |
| M1 | 0.21 | 22.1 | 0.18 | 340 |
| CL | 0.28 | 19.1 | 0.24 | 191 |
| CC | 0.18 | 24.1 | 0.15 | 357 |
| EL | 0.11 | 21.6 | 0.09 | 371 |
| EC | 0.19 | 17.7 | 0.08 | 350 |
| ER | 0.05 | 18.8 | 0.13 | 291 |
| All | **0.18** | **19.8** | **0.10** | **323** |
| LSMs | 0.21 | 18.7 | 0.05 | 287 |
| Coupled | 0.18 | 19.3 | 0.10 | 352 |
| Reanalyses | 0.15 | 21.4 | 0.16 | 328 |

788
789

790 **Figure Captions:**

791  Figure 1: Location of the FLUXNET2015 Tier-1 sites used in this study. Triangles indicate

792    no upward shortwave radiation measurements available to estimate surface albedo,

793    pluses mean no Bowen ratio corrected surface heat fluxes provided, exes indicate neither

794    albedo nor corrected heat fluxes are available, circles have both. Color of the symbol

795    indicates the length of data series available.

796  Figure 2: Scatter of annual total precipitation measurements at FLUXNET2015 sites

797    (abscissa) to estimates (ordinate) from gridded observationally-based precipitation

798    analyses (top two rows) or reanalyses constrained by data assimilation (bottom row)

799    using the value from the grid box containing the FLUXNET2015 site location (unless data

800    are missing or indicated to be an all-ocean grid box). Dash-dotted diagonal grey line

801    indicates X=Y. Colors indicate years of available data from each FLUXNET2015 site,

802    whiskers span range of annual totals from FLUXNET2015 (horizontal) or gridded

803    estimates (vertical) for years where data sets overlap. Purple line is the best-fit linear

804    regression of Y on X. Statistics are explained in the text.

805  Figure 3: Taylor diagram of annual mean surface radiation terms for the 12 indicated

806    models verified against FLUXNET2015 sites for downward solar radiation (black),

807    downward longwave radiation (red), upward shortwave radiation (blue) and net

808    radiation (green). Dot colors indicate mean bias and size shows percentage of stations

809    where the range of the annual totals from the model overlaps the span from

810    FLUXNET2015 sites (also presented in tabular form in the upper right).

811  Figure 4: As in Fig. 3 for surface albedo; annual mean (black) and boreal summer (JJA)

812    mean (red).

813   Figure 5: Propagation of errors estimated from their rank correlations among precipitation

814     (P), height of the lifting condensation level (LCL), evaporative fraction (EF), sensible and

815     latent heat flux (SH & LH), surface albedo ($\alpha$) and net radiation ($R_{Net}$) across

816     FLUXNET2015 stations. Ratios show the number of models out of 11 (correlations

817     involving $\alpha$) or 12 (other variables) with p-values below 0.10; p-value shown is based on

818     the average of correlations across all models. Widths of arrows follow significance of

819     correlations and no arrows are drawn where p-values are large. The wide double arrows

820     between EF and heat fluxes denote p-values $< 10^{-12}$.

821   Figure 6: As in Fig. 3 for the magnitude of the annual cycle (first harmonic calculated from

822     monthly means) of sensible heat flux (orange), latent heat flux (cyan) and net radiation at

823     the surface (green).

824   Figure 7: As in Fig. 6 for phase of the annual cycle of sensible heat flux (orange) and latent

825     heat flux (cyan) and net radiation at the surface (green).

826   Figure 8: As in Fig. 6 for the magnitude (brown) and phase (purple) of the annual cycle of

827     surface albedo.

828   Figure 9: Distribution of coupling indices for the terrestrial (x-axis) and atmospheric (y-

829     axis) legs for the warmest consecutive 3 months of the annual cycle for FLUXNET2015

830     sites (white dots; identical in each panel) and for each model as indicated. Colors of dots

831     indicate in which quadrant that FLUXNET2015 site lies: red = both indices positive;

832     green = terrestrial positive, atmospheric negative; blue = atmospheric positive,

833     terrestrial negative; grey = both negative. The white circle indicates the centroid of all

834     FLUXNET2015 stations that are in that quadrant, connected by a colored dotted line to a

835     colored circle that is the centroid of the same stations' corresponding grid boxes as

836     simulated by the model. Numbers in the corners of each quadrant show the number of

837      points in that quadrant according to the model and FLUXNET2015 data, separated by a

838      colon, and the percentage of the FLUXNET2015 sites within that quadrant that the model

839      placed in the same quadrant. The percentage in red at the upper right of each panel is the

840      overall percentage of sites where model and FLUXNET2015 agree on the quadrant.
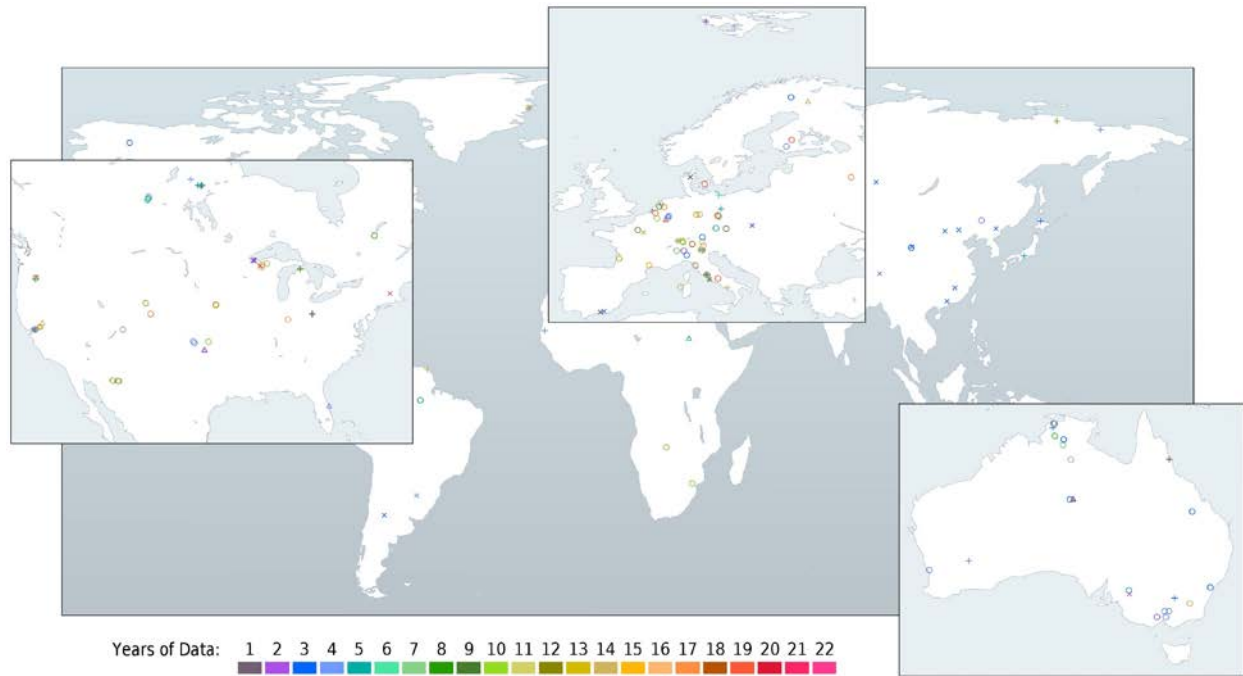
841

842

843 Figure 1: Location of the FLUXNET2015 Tier-1 sites used in this study. Triangles indicate
844 no upward shortwave radiation measurements available to estimate surface albedo, pluses
845 mean no Bowen ratio corrected surface heat fluxes provided, exes indicate neither albedo
846 nor corrected heat fluxes are available, circles have both. Color of the symbol indicates the
847 length of data series available.

848

Figure 2: Scatter of annual total precipitation measurements at FLUXNET2015 sites (abscissa) to estimates (ordinate) from gridded observationally-based precipitation analyses (top two rows) or reanalyses constrained by data assimilation (bottom row) using the value from the grid box containing the FLUXNET2015 site location (unless data are missing or indicated to be an all-ocean grid box). Dash-dotted diagonal grey line indicates X=Y. Colors indicate years of available data from each FLUXNET2015 site, whiskers span range of annual totals from FLUXNET2015 (horizontal) or gridded estimates (vertical) for years where data sets overlap. Purple line is the best-fit linear regression of Y on X. Statistics are explained in the text.

860

Figure 3: Taylor diagram of annual mean surface radiation terms for the 12 indicated
models verified against FLUXNET2015 sites for downward solar radiation (black),
downward longwave radiation (red), upward shortwave radiation (blue) and net radiation
(green). Dot colors indicate mean bias and size shows percentage of stations where the
range of the annual totals from the model overlaps the span from FLUXNET2015 sites (also
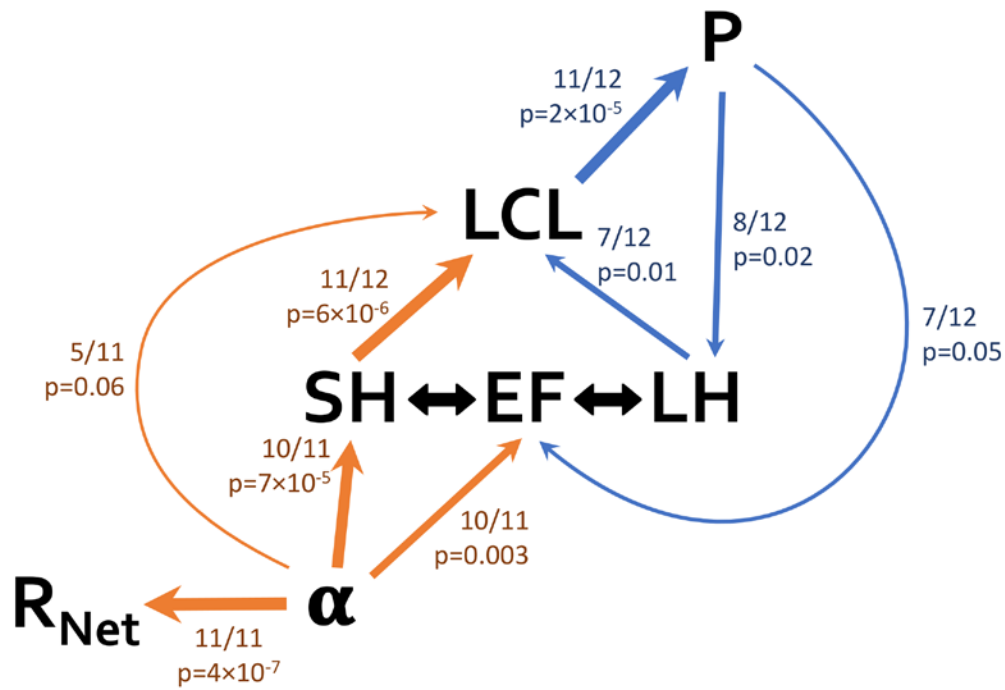presented in tabular form in the upper right).

867

868

Figure 4: As in Fig. 3 for surface albedo; annual mean (black) and boreal summer (JJA) mean (red).

Figure 5: Propagation of errors estimated from their rank correlations among precipitation (P), height of the lifting condensation level (LCL), evaporative fraction (EF), sensible and latent heat flux (SH & LH), surface albedo ($\alpha$) and net radiation ($R_{Net}$) across FLUXNET2015 stations. Ratios show the number of models out of 11 (correlations involving $\alpha$) or 12 (other variables) with p-values below 0.10; p-value shown is based on the average of correlations across all models. Widths of arrows follow significance of correlations and no arrows are drawn where p-values are large. The wide double arrows between EF and heat fluxes denote p-values $< 10^{-12}$.
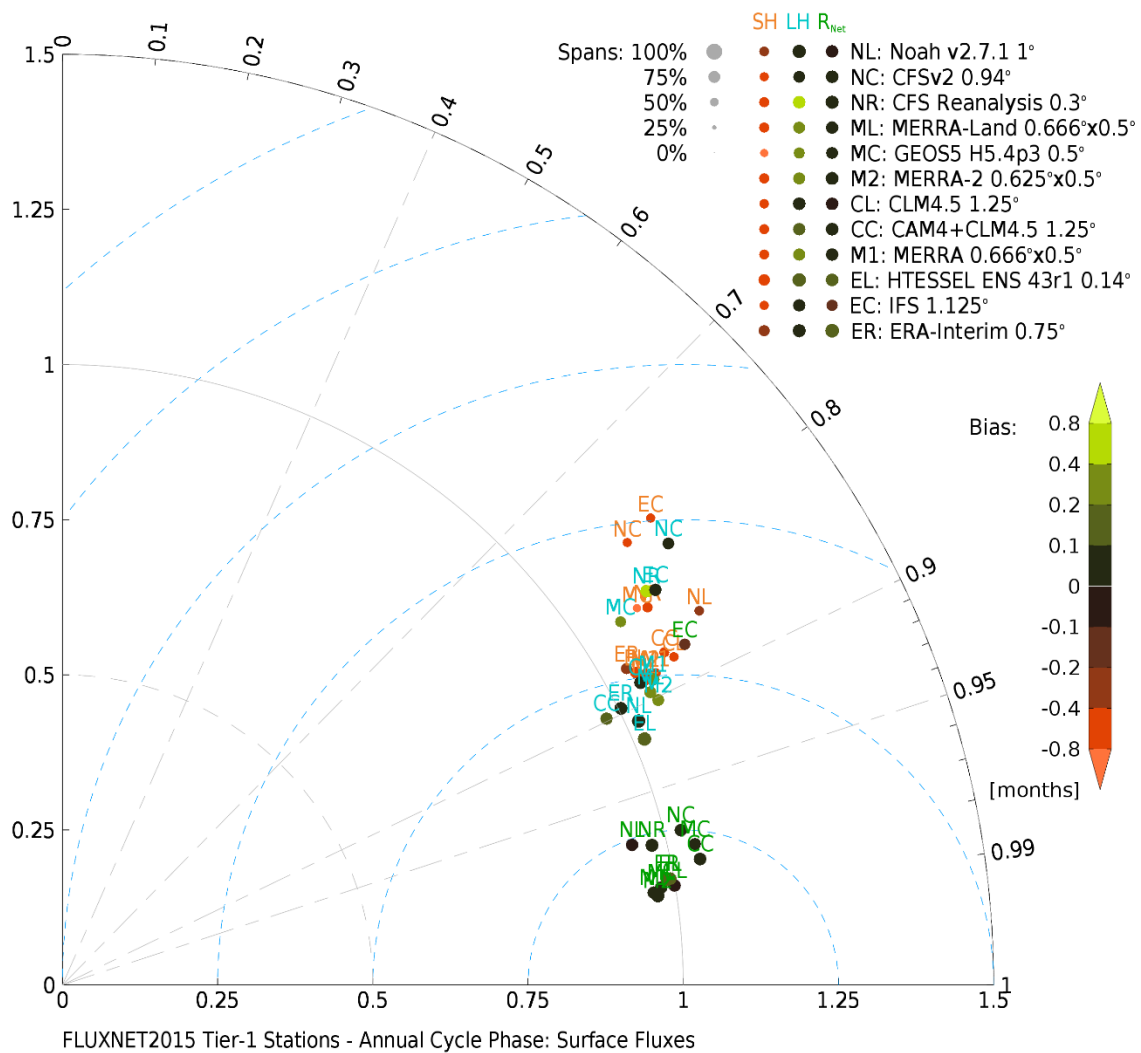
Figure 6: As in Fig. 3 for the magnitude of the annual cycle (first harmonic calculated from monthly means) of sensible heat flux (orange), latent heat flux (cyan) and net radiation at the surface (green).
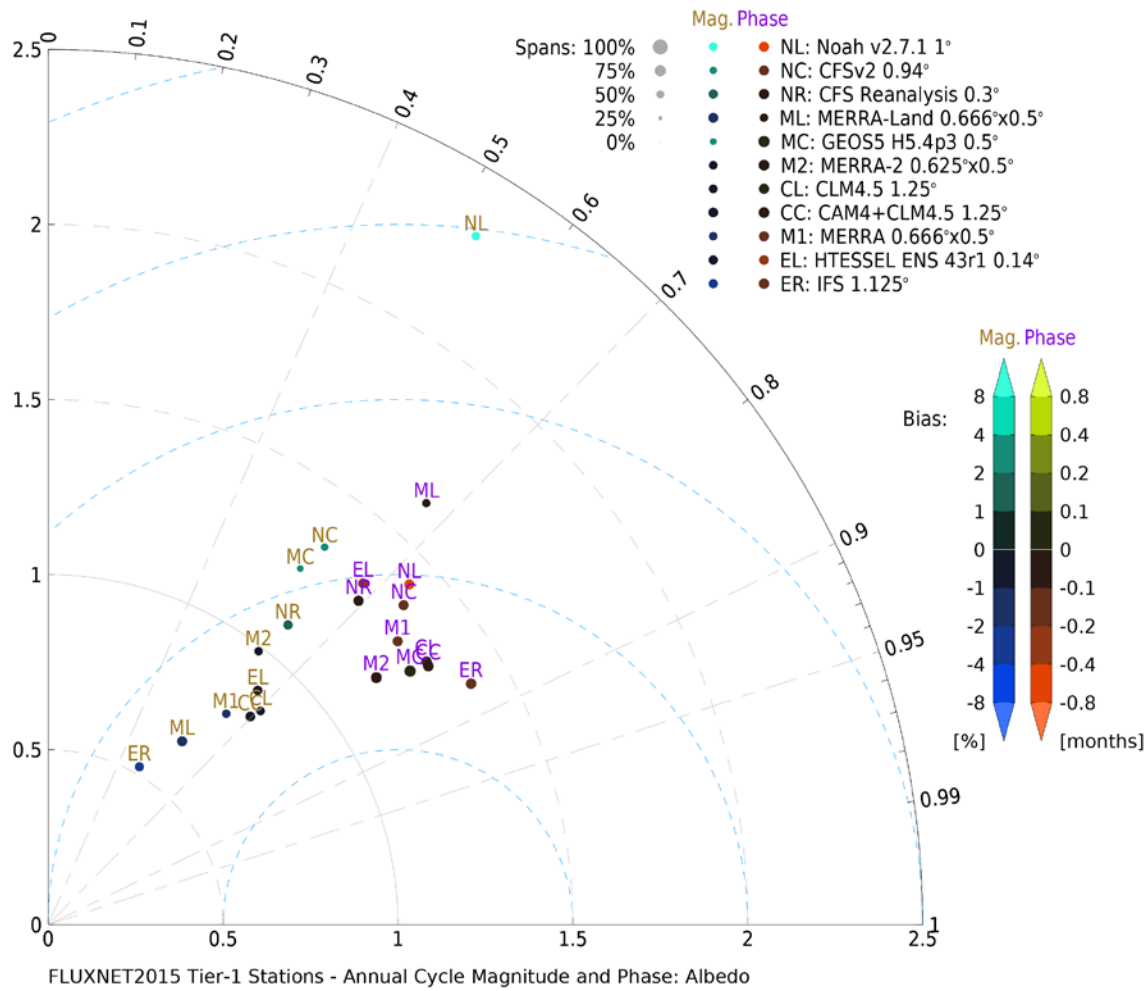
888

Figure 7: As in Fig. 6 for phase of the annual cycle of sensible heat flux (orange) latent heat

flux (cyan), and net radiation at the surface (green).

891

Figure 8: As in Fig. 6 for the magnitude (brown) and phase (purple) of the annual cycle of surface albedo.
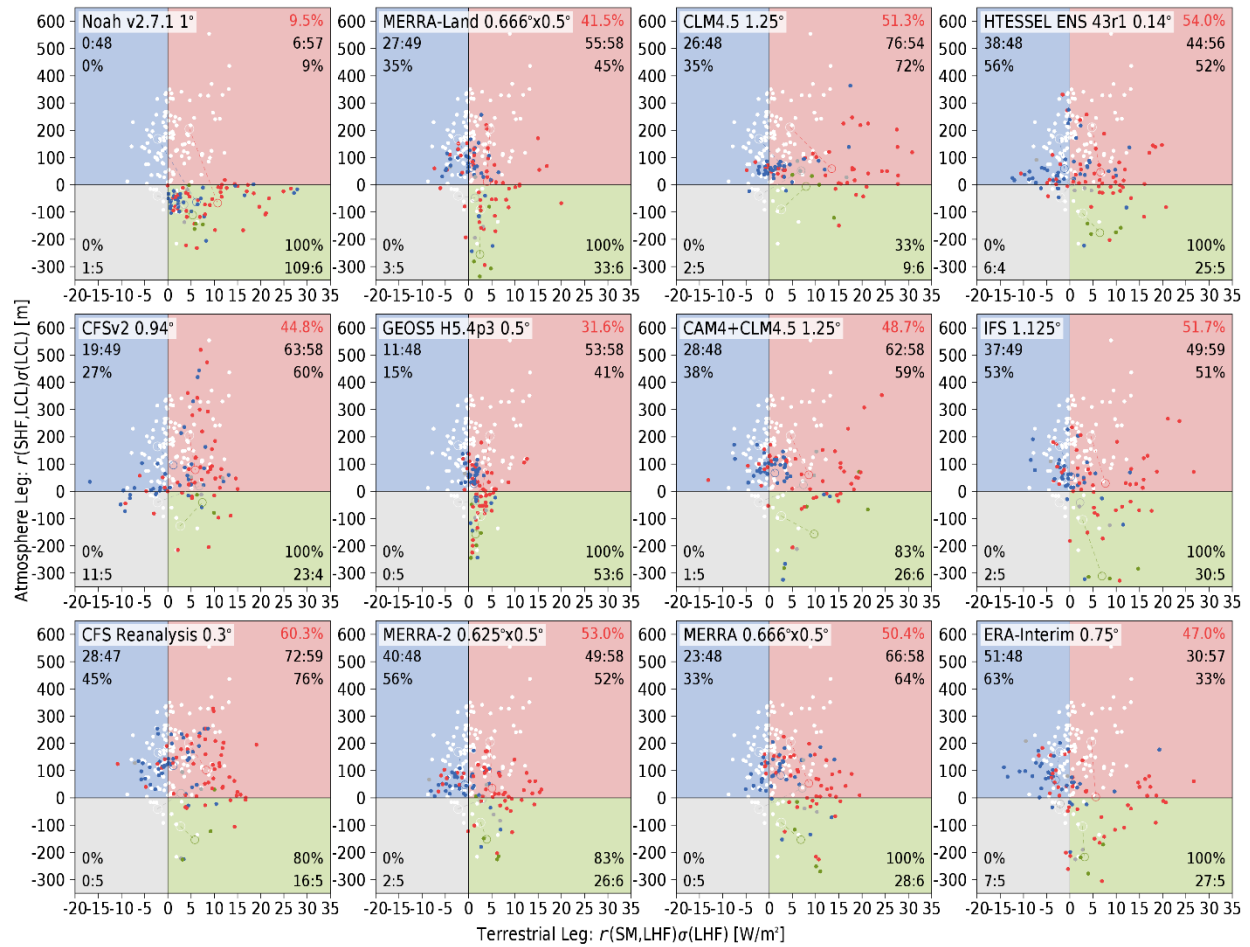
Figure 9: Distribution of coupling indices for the terrestrial (x-axis) and atmospheric (y-axis) legs for the warmest consecutive 3 months of the annual cycle for FLUXNET2015 sites (white dots; identical in each panel) and for each model as indicated. Colors of dots indicate in which quadrant that FLUXNET2015 site lies: red = both indices positive; green = terrestrial positive, atmospheric negative; blue = atmospheric positive, terrestrial negative; grey = both negative. The white circle indicates the centroid of all FLUXNET2015 stations that are in that quadrant, connected by a colored dotted line to a colored circle that is the centroid of the same stations' corresponding grid boxes as simulated by the model. Numbers in the corners of each quadrant show the number of points in that quadrant according to the model and FLUXNET2015 data, separated by a colon, and the percentage of the FLUXNET2015 sites within that quadrant that the model placed in the same quadrant. The percentage in red at the upper right of each panel is the overall percentage of sites where model and FLUXNET2015 agree on the quadrant.